

# DATA STORAGE SYSTEM

Publication number: JP9198308 (A)

Publication date: 1997-07-31

Inventor(s): FUJIBAYASHI AKIRA; TAKAMOTO YOSHIFUMI

Applicant(s): HITACHI LTD

Classification: G06F12/08; G06F3/06; G06F12/08; G06F3/06; (IPC1-7): G06F12/08; G06F3/06; G06F12/08

- European:

Application number: JP19960023202 19960117

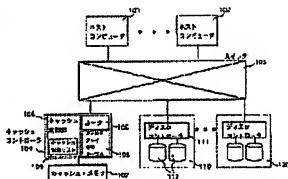
Priority number(s): JP19960023202 19960117

Also published as:

JP3776496 (B2)

## Abstract of JP 9198308 (A)

**PROBLEM TO BE SOLVED:** To effectively use a cache memory and to improve data access performance in a data storage system using a switch. **SOLUTION:** Plural host computers 101 and 102 are connected with secondary storage devices 110 and 120 through the switch 103. A cache memory 107 provided with a cache controller 104 common to the secondary storage devices 110 and 120 is switch-connected to the secondary storage devices in parallel. At the time of accessing data, the host computers 101 and 102 transmit data to the cache memory 107 through the switch. When a cache error occurs in the cache memory 107, a disk array management table 108 is referred to, a port to which the secondary storage device corresponding to a logic volume in a packet is connected is specified and the secondary storage device is accessed. Data transmitted from the secondary storage device is stored in the cache memory 107 and data is transferred to the host computers 101 and 102 through the switch 103.



特開平9-198308

(43) 公開日 平成9年(1997)7月31日

(51) Int.Cl.*	識別記号	序内整理番号	F I	技術表示箇所
G 0 6 F 12/08		7623-5B	G 0 6 F 12/08	H
		7623-5B		P
	3 2 0	7623-5B		3 2 0
3/08	3 0 1		3/08	3 0 1 M
	3 0 2			3 0 2 A

審査請求 未請求 請求項の数 5 F D (全 10 頁)

(21) 出願番号 特願平8-23202

(22) 出願日 平成8年(1996)1月17日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 藤林 昭

東京都国分寺市東荏ケ窪1丁目280番地

株式会社日立製作所中央研究所内

(72) 発明者 高本 良史

東京都国分寺市東荏ケ窪1丁目280番地

株式会社日立製作所中央研究所内

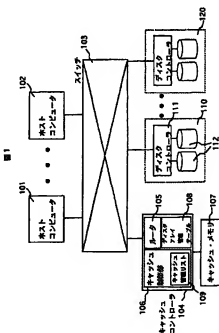
(74) 代理人 弁理士 菅岡 茂 (外1名)

## (54) 【発明の名称】 データ記憶システム

## (57) 【要約】

【課題】 スイッチを用いたデータ記憶システムにおいて、キャッシュメモリを有効に活用してデータアクセス性能を向上することにある。

【解決手段】 複数のホストコンピュータ101、102と2次記憶装置110、120はスイッチ103を介して接続され、2次記憶装置に共通のキャッシュコントローラ104を備えるキャッシュメモリ107が2次記憶装置と並列にスイッチに接続される。ホストコンピュータはデータをアクセスするときパケットをスイッチを介してキャッシュメモリに送出する。キャッシュメモリでキャッシュミスが生じたときは、ディスクアレイ管理テーブル107を参照してパケット内の論理ボリュームに対応する2次記憶装置が接続されているポートを特定して2次記憶装置をアクセスし、2次記憶装置から送られてくるデータをキャッシュメモリに格納し、データをスイッチを介してホストコンピュータに転送する。



【特許請求の範囲】

【請求項1】 複数のホストコンピュータと、複数の2次記憶装置とから構成され、該複数のホストコンピュータと該複数の2次記憶装置との間をスイッチにより接続するデータ記憶システムにおいて、該スイッチに該複数の2次記憶装置と並列に独立したキャッシュメモリを接続し、該キャッシュメモリは該複数の2次記憶装置の該スイッチを介したキャッシュメモリであることを特徴とするデータ記憶システム。

【請求項2】 請求項1記載のデータ記憶システムにおいて、

前記キャッシュメモリは、前記ホストコンピュータのデータ転送要求に対して、要求されたデータを検索し、検索の結果該データがキャッシュメモリ内にない場合には、前記スイッチを介して該データが格納されている2次記憶装置から該データを取り込み、該キャッシュメモリに格納した後にホストコンピュータに該データを転送することを特徴とするデータ記憶システム。

【請求項3】 請求項2記載のデータ記憶システムにおいて、

前記キャッシュメモリはキャッシュコントローラを備え、該キャッシュコントローラは、複数の各2次記憶装置が接続される前記スイッチのポート番号とホストコンピュータの認識している論理ボリュームを対応させたディスクアレイ管理テーブルを有し、前記検索の結果前記データがキャッシュメモリ内にない場合、該ディスクアレイ管理テーブルを参照して2次記憶装置をアクセスすることを特徴とするデータ記憶システム。

【請求項4】 請求項3記載のデータ記憶システムにおいて、

前記各ホストコンピュータのオペレーションシステムに前記ディスクアレイ管理テーブルを設け、該各ホストコンピュータは、前記複数の2次記憶装置のいずれかに直接アクセスするとき、該ディスクアレイ管理テーブルを参照して2次記憶装置をアクセスすることを特徴とするデータ記憶システム。

【請求項5】 請求項3記載のデータ記憶システムにおいて、

前記スイッチは、前記各ホストコンピュータが前記キャッシュメモリを使用するか直接前記2次記憶装置を使用するかを示すテーブルと直接前記2次記憶装置を使用する場合に用いられる前記ディスクアレイ管理テーブルと同様のテーブルからなるホスト管理テーブルを有し、前記ホストコンピュータからのアクセス要求に応じて前記ホスト管理テーブルを参照して前記キャッシュメモリまたは前記2次記憶装置に前記ホストコンピュータからのアクセスデータを転送することを特徴とするデータ記憶システム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、複数のホストコンピュータと複数の2次記憶装置をスイッチを利用して接続する構成のデータ記憶システムに関する。

【0002】

【従来の技術】 一般的なデータ処理システムは、ホストコンピュータと2次記憶装置から構成されている。2次記憶装置として使用されるのは主に磁気ディスク装置である。ここで発明者のいう磁気ディスク装置は単体のディスクドライブまたはディスクアレイを意味する。本発明の利用分野として挙げたスイッチを利用したデータ記憶システムの一例が特開平7-44322号において記述されている。

【0003】 スwitchを用いてホストコンピュータと磁気ディスク装置を接続する構成のデータ記憶システムでは、ホストコンピュータおよび磁気ディスク装置の台数の変更にも柔軟に対応することが可能でスケラビリティに優れる。このスイッチを利用したホストコンピュータと磁気ディスク装置間のデータ転送はパケット交換によって実現される。これを簡単に説明すると、ホストコンピュータはコマンドパケット（例えば、ポート番号1番；セクタ番号2；データサイズ512 Byte；Readのような構成）を送信し、そのコマンドパケットをスイッチが解析してホストコンピュータと要求されたポート間の接続が確立したのちに、ホストコンピュータとそのポートに接続している磁気ディスク装置間でデータ転送を開始すると言う手順になる。また、従来のバススイッチと異なり、ある1組のホストコンピュータとの磁気ディスク装置の間でデータ転送が行われていると、データバスは占有されず、同時に他のホストコンピュータと磁気ディスク装置の間でもデータ転送が行える。ここで、

データの読出しに関して、磁気ディスク装置では、ホストコンピュータから頻繁にアクセス要求のあるデータに対して、その都度読出し動作を行うと、メカニカルな動作を伴うためデータ転送に時間がかかるという問題がある。この問題を解決するために、磁気ディスク装置は通常キャッシュメモリを備えており、頻繁にアクセス要求のあるデータは、そのキャッシュメモリに格納することでそのデータの読出しに対するメカニカルな動作を省略し性能の向上を図っている。しかし、従来の例として挙げたシステム構成では、従来のバススイッチを本発明の用いているスイッチに置き換え、ホストコンピュータとのインターフェース以下の部分はアレイコントローラが磁気ディスク装置とホストコンピュータ間のデータ転送を制御している従来の一般的構成である。この様な構成の場合には、接続された磁気ディスク装置の各々にキャッシュメモリが分散されて配置されており、それぞれのキャッシュが有効に利用されていない。

【0004】

【発明が解決しようとする課題】 従来例として挙げたデータ記憶システムにおいて、接続している複数の磁気ディ

スク装置のそれぞれにキャッシュメモリが備わっている場合、性能向上のためのキャッシュメモリが、接続している磁気ディスク装置それぞれに分散して配置されることになり、それぞれのキャッシュを有効に利用できない。例えば、ある一つの磁気ディスク装置に格納されているデータに頻繁にアクセス要求が送られてくるような場合には、他の磁気ディスク装置に備わっているキャッシュメモリはほとんど利用されないことになり無駄になってしまう。

【0005】本発明の目的は、スイッチを用いたデータ記憶システムにおいて、キャッシュメモリを有効に活用してデータアクセス性能を向上することにある。

【0006】

【課題を解決するための手段】上記目的を達成するため、本発明は、複数のホストコンピュータと、複数の2次記憶装置とから構成され、該複数のホストコンピュータと該複数の2次記憶装置との間をスイッチにより接続するデータ記憶システムにおいて、該スイッチに該複数の2次記憶装置と並列に独立したキャッシュメモリを接続し、該キャッシュメモリは該複数の2次記憶装置の該スイッチを介したキャッシュメモリであるようにしている。また、前記キャッシュメモリは、前記ホストコンピュータのデータ転送要求に対して、要求されたデータを検索し、検索の結果該データがキャッシュメモリ内にある場合には、前記スイッチを介して該データが格納されている2次記憶装置から該データを取り込み、該キャッシュメモリに格納した後にホストコンピュータに該データを転送するようにしている。前記キャッシュメモリはキャッシュコントローラを備え、該キャッシュコントローラは、複数の各2次記憶装置が接続される前記スイッチのポート番号とホストコンピュータの認識した物理ボリュームを対応させたディスクアレイ管理テーブルを有し、前記検索の結果前記データがキャッシュメモリ内にある場合、該ディスクアレイ管理テーブルを参照して2次記憶装置をアクセスするようにしている。また、前記各ホストコンピュータのオペレーションシステムに前記ディスクアレイ管理テーブルを設け、該各ホストコンピュータは、前記複数の2次記憶装置のいずれかに直接アクセスするとき、該ディスクアレイ管理テーブルを参照して2次記憶装置をアクセスするようにしている。また、前記スイッチは、前記各ホストコンピュータが前記キャッシュメモリを使用するか直接前記2次記憶装置を使用するかを示すテーブルと直接前記2次記憶装置を使用する場合に用いられる前記ディスクアレイ管理テーブルと同様のテーブルからなるホスト管理テーブルを有し、前記ホストコンピュータからのアクセス要求に応じて前記ホスト管理テーブルを参照して前記キャッシュメモリまたは前記2次記憶装置に前記ホストコンピュータからのアクセスデータを転送するようにしている。

【0007】

【実施例】本発明の提供するデータ記憶システムを以下に図面を示し実施例を参照して詳細に説明する。図1は、本発明によるデータ記憶システムの構成をブロック図で示したものである。101、102はホストコンピュータであり、103はスイッチであり、104はキャッシュコントローラであり、ルータ105とキャッシュ制御部106から構成される。107はキャッシュメモリである。108はキャッシュコントローラ104の持つディスクアレイ管理テーブルである。109はキャッシュ制御部106がキャッシュ検索時に参照するキャッシュ管理リストである。110、120は磁気ディスク装置であり、この図1中では一例としてディスクアレイ装置としており、ディスクコントローラ111、複数のディスク装置112から構成されている。磁気ディスク装置120も同様の構成である。この図に示すように本発明はキャッシュメモリ107を、複数の磁気ディスク装置110、120と共にスイッチ103に並列に接続することで、キャッシュメモリ107は前記磁気ディスク装置110、120及びホストコンピュータ101、102と、同様の1つのデバイスのごとくスイッチにより容易にアクセスできる構成となる。そして、磁気ディスク装置110、120間にはキャッシュメモリを設けず、キャッシュメモリ107を磁気ディスク装置110、120が共用する。

【0008】ここで、スイッチ103の内部構造の一例を図7に示す。701はデータのシリアル/パラレル変換を行う部分(S/P)である。702はスイッチ制御装置である。703はスイッチ機構である。ここで、スイッチ103の動作の一例を簡単に説明する。図7の中に示すような構成のポート704がポート1に接続されているホストコンピュータから発行される。パケットは送信先ポート番号とデータ部からなる。データ部の内容は705に示すように、コマンド、論理ボリューム(Lvol#)、ブロック番号(BLK#)、ポート番号(Host#)からなる。送信されてきたパケットの送信先ポートをみて、スイッチ制御装置702がそのポートとの接続を確認するように制御線を通じて信号を送出する。送信先ポートとの接続が確立したのち、送信先に向けてパケットまたはデータが送られる。

【0009】次に上記システム構成においてキャッシュメモリ107を有効利用するために用いるデータ転送方法について述べる。図1のキャッシュ制御部106で行うホストコンピュータからのデータ読出し要求に対するデータ転送のフローチャートを図2に示す。まず、ホストコンピュータはデータ読出し要求をキャッシュメモリの接続ポートに対して送る。ステップ201において、要求されたデータをキャッシュ制御部108において、送信されたパケット内の論理ボリュームLvolとブロック番号BLKに基づいてキャッシュ管理リスト109を参照しキャッシュメモリ107内にデータが存在す

るかどうか検索する。次に判断ステップ202において、要求されたデータがキャッシュメモリ107内に存在する(キャッシュヒット)場合には、ステップ205において、要求されたデータをホストコンピュータに転送する。要求されたデータが存在しない(キャッシュミス)場合には、ステップ203において、ルータに制御を移す。続いてステップ204において、ルータの制御で磁気ディスク装置から送られて来たデータをキャッシュ制御部106がキャッシュメモリ107に格納し、キャッシュ管理リストを更新する。この時、キャッシュメモリの容量一杯までデータがすでに格納されている場合は、データの追い出しが必要となる。現在追い出しの手法は様々な方法が利用されているが、ここでは最も使用頻度が少なく最も古いデータを追い出しの候補に選ぶ方法を使用することとするが、その他の方法を用いたとしても本発明の実施にはなんら問題は無い。そして、ステップ205において、ホストコンピュータに要求されたデータを転送する。

【0010】ここで、図2のステップ201においてキャッシュ制御部が参照するキャッシュ管理リストの概略を図9に示す。キャッシュコントローラ内のルータが受信したパケット内の論理ボリューム番号とブロック番号がキャッシュ制御部に渡される。キャッシュメモリ内のデータはキャッシュ管理リストによりブロック単位に管理されている。現在使用中(キャッシュメモリ内にデータが存在する)論理ボリューム番号とブロック番号のリストと未使用のリストを持っており、渡された論理ボリューム番号とブロック番号が使用リスト中に存在するかどうか検索する。

【0011】図3は前記ルータに制御を移した後の処理をフローチャートで示したものである。ここで、ルータの機能は、簡略に説明すると、ホストコンピュータから転送されたパケット内のディスクコマンドを解析し、その結果に基づいて所定の磁気ディスク装置を選択しコマンドやデータをルーティングすることである。図3中のステップ301において、ルータは送信されてきたパケット内の論理ボリュームをみて、ディスクアレイ管理テーブル(以下、図4において説明する。)を参照し、その論理ボリュームに対応する磁気ディスク装置が接続されているポートを特定する。ステップ302において、特定したポートに対しパケットを送出する。ステップ303において、このパケットを受けた磁気ディスク装置から転送されてくるデータを受け取る。

【0012】これまで述べてきた実施例の通り、ホストコンピュータへのデータ転送時は常時キャッシュメモリ107を使用することで、例えば磁気ディスク装置101に格納されているデータに頻繁にアクセスがあるような場合でも、他の磁気ディスク装置に格納されているデータもキャッシュメモリ107に格納しておくことで、すべての磁気ディスク装置に対してキャッシュメモ

リ107が有効利用される。

【0013】図4は前記ディスクアレイ管理テーブルである。ホストコンピュータのオペレーティングシステム(以下、OSと記述する)は、この論理ボリューム(Logical)とブロック(Block)によりデータを指定してアクセス要求を発行する。実際のデータは複数の磁気ディスク装置内に格納されているので、その磁気ディスク装置の接続されているポートとの対応を取るためにこのディスクアレイ管理テーブルを用いる。このディスクアレイ管理テーブルをキャッシュコントローラ104が持ち、前記実施例のようにホストコンピュータへのデータ転送時に常時キャッシュメモリ107を利用することで、ホストコンピュータのOSは論理ボリュームを用いてデータアクセス要求を発行することになり、複数の磁気ディスク装置が接続されていることとその内部のデータ配列の構成をユーザーに意識させず、単一の磁気ディスク装置の使用環境を提供する。

【0014】図5では本発明の提供するシステムにおいて、実際のデータ転送の一例を簡略に示している。ここでは、ホストコンピュータ101からデータ1への読出し要求が発行され、ホストコンピュータ102からデータ2への読出し要求が発行されている場合を考える。データ1は前記キャッシュメモリ内に現在格納されており、データ2はこのキャッシュメモリ内に無く磁気ディスク装置120に格納されている。ホストコンピュータ101からのデータ1の読出し要求を図5中に示すような構成のパケット501として送信する。図2で示したフローチャートに従い、先ず最初に上記キャッシュメモリ107内を検索される。データ1はキャッシュヒットするので、キャッシュメモリ107からホストコンピュータ101に転送される。同様にホストコンピュータ102はパケット502を送信して、上記キャッシュメモリ107内を検索されるがキャッシュミスとなり、キャッシュコントローラ内のルータが、図4で説明したディスクアレイ管理テーブルを参照してデータ2の格納位置を磁気ディスク装置120との接続されているポートと特定し、これに対してデータ2の読出し要求のパケット503を発行して、データ2を取り込み、上記キャッシュメモリ107内に格納してから、ホストコンピュータ102に転送する。

【0015】これまで説明した実施例では接続するすべてのホストコンピュータは本発明の提供するキャッシュメモリにアクセスするシステムになっている。しかし、ホストコンピュータによっては、使用するアプリケーションの性質により、本発明の提供するキャッシュメモリより直接磁気ディスク装置にアクセスする方が性能が良いという場合がある。これに対して、以下に説明する実施例で対応する。一つの実施例は、ディスクアレイ管理テーブルをホストコンピュータのOSにも持たせることである。これにより、ホストコンピュータには論理

ボリュームと実際にデータの格納されている磁気ディスク装置が接続されているポートの対応が分かっているもので、発行するバケットの送信先ポートにキャッシュメモリまたは磁気ディスク装置の接続されているポートを指定することで、本発明のキャッシュメモリの使用、不使用が選択可能になる（概略図を図6に示す。）

もう一つの実施例は前記スイッチにホストコンピュータのアクセス要求の管理を行う構成管理テーブルを持たせることである。図8にその概略図を示す。このホスト管理テーブル801は、スイッチ103内のスイッチ制御装置702に接続され、送信されてきたバケット内のホスト番号とホスト管理テーブル801を照合して、キャッシュメモリ107を使用するホストコンピュータと使用しないホストコンピュータを判断し、ホストコンピュータの発行する転送要求先をキャッシュメモリ107か、要求するデータの格納されている磁気ディスク装置かに振り分ける。例えば、ホストコンピュータの発行するアクセス要求をホスト管理テーブル801を参照してキャッシュメモリを使用しないホストである場合、ホスト管理テーブルはディスクアレイ管理テーブルと同様の

【0016】

【発明の効果】本発明により、スイッチを用いたデータ記憶システムにおいて、キャッシュメモリを有効に活用してデータアクセス性能が向上する。さらに、ホストコンピュータの使用するアプリケーションによってキャッシュメモリを使用するか、キャッシュメモリを不使用として直接磁気ディスク装置にアクセスするかの選択性を有する効率のよいデータ記憶システムを可能とする。

【図面の簡単な説明】

【図1】本発明によるキャッシュメモリを備えるスイッチを利用したデータ記憶システムの構成を示すブロック図である。

【図2】ホストコンピュータからの読出し要求に対するキャッシュ制御部の処理のフローチャートを示す図である。

【図3】ルータがデータをキャッシュメモリに格納する場合の処理のフローチャートを示す図である。

【図4】ディスクアレイ管理テーブルの一例を示す図である。

【図5】本発明における実際のデータの流れの一例を示した概略図である。

【図6】ディスクアレイ管理テーブルをホストコンピュータのOSに持たせた場合の本発明のデータ記憶システムの構成を示すブロック図である。

【図7】スイッチの構造の概略図である。

【図8】スイッチがホスト管理テーブルを備える場合の本発明のデータ記憶システムの構成を示すブロック図である。

【図9】キャッシュ管理リストの概略図である。

【符号の説明】

101、102 ホストコンピュータ

103 スイッチ

104 キャッシュコントローラ

105 ルータ

106 キャッシュ制御部

107 キャッシュメモリ

108 ディスクアレイ管理テーブル

110、120 磁気ディスク装置

111 ディスクコントローラ

112 ディスク装置

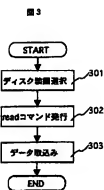
701 S/P

702 スイッチ制御装置

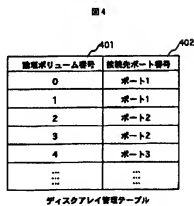
703 スイッチ機構

801 ホスト管理テーブル

【圖3】



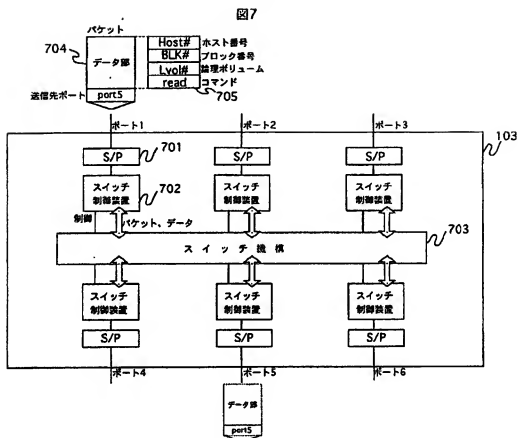
【圖4】



[illegible]

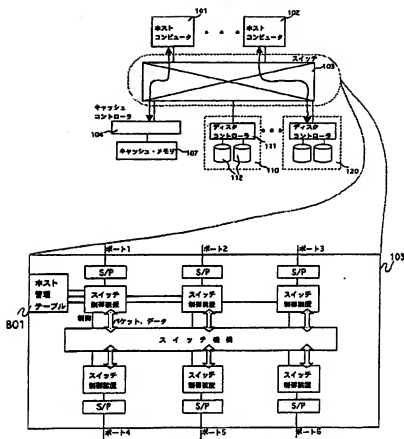


〔図7〕



〔図8〕

図8



ホスト管理テーブル

host0	use	Lvol0	port4
host1	no use	Lvol1	port4
host2	use	Lvol2	port5
host3	no use	Lvol3	port5
host4	use	Lvol4	port6
host5	use	Lvol5	port6
⋮	⋮	⋮	⋮

【図9】

図9

